

Feb 21, 2020

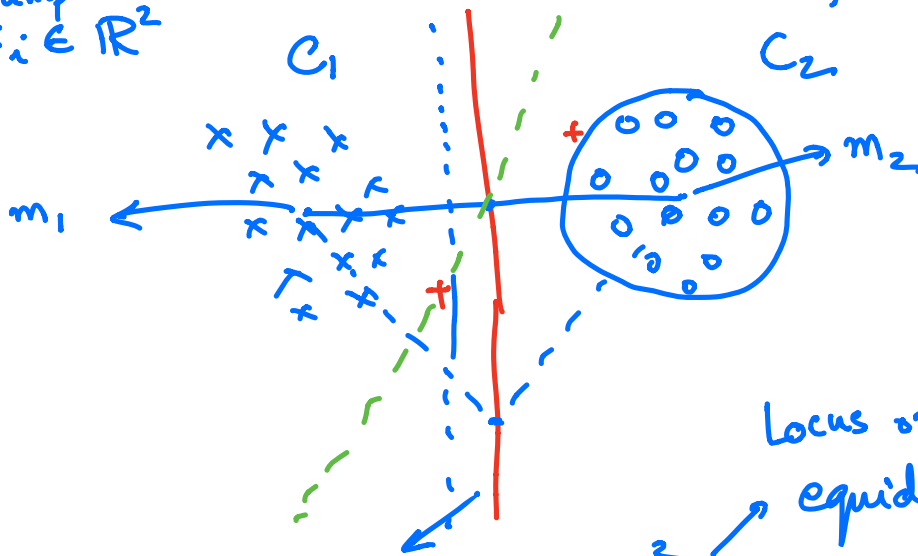
# Classification

Training set:  $(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)$

Goal: Learn  $f$  that predicts class

label of new (test) point  $x$ .

Example:  
 $x_i \in \mathbb{R}^2$



$$\|x - m_1\|_2^2 = \|x - m_2\|_2^2 \quad \text{Locus of points equidistant to } m_1 \text{ \& } m_2$$

$$\cancel{x^T x} - 2x^T m_1 + \|m_1\|_2^2 = \cancel{x^T x} - 2x^T m_2 + \|m_2\|_2^2$$

$$y(x) = (m_2 - m_1)^T x + \frac{1}{2} (\|m_1\|_2^2 - \|m_2\|_2^2) = 0$$

$$y(x) > 0, x \in C_2$$

$$y(x) < 0, x \in C_1$$

↓  
 Example of a hyperplane (linear variety of (d-1) dimensions when d is dimensionality of data)

H:  $w^T x + w_0 = 0$  - Hyperplane

If  $x$  lies on the hyperplane,

$$w^T x + w_0 = 0$$

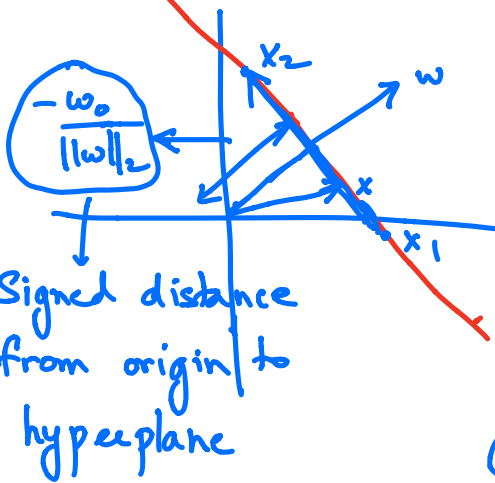
$$\left(\frac{w}{\|w\|_2}\right)^T x = -\frac{w_0}{\|w\|_2}$$

$$w^T x_1 + w_0 = 0$$

$$w^T x_2 + w_0 = 0$$

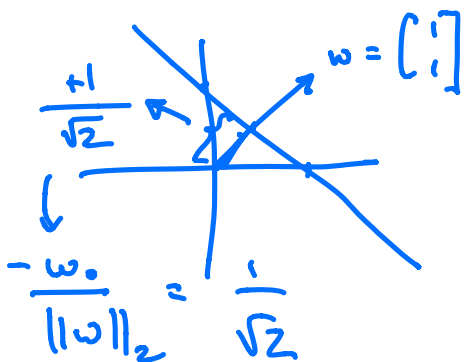
(subtract)  $w^T (x_1 - x_2) = 0$

$w$  is normal to (the points on) hyperplane



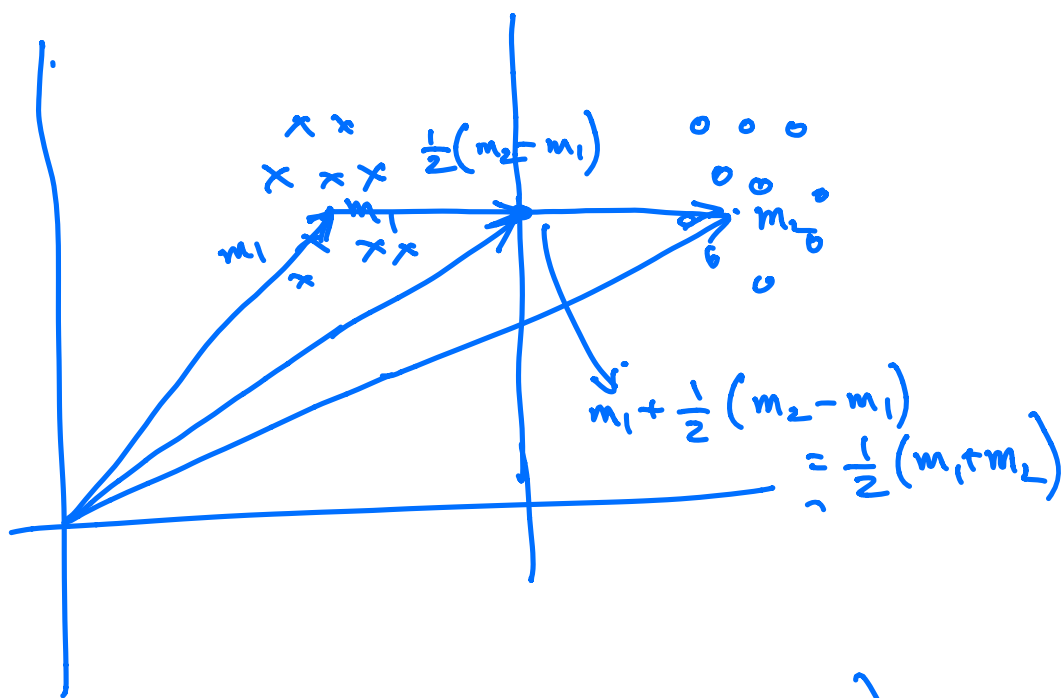
Signed distance from origin to hyperplane

Example:  $x_1 + x_2 = 1$



$$\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - 1 = 0$$

$$w = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, w_0 = -1$$



$$y(x) = (m_2 - m_1)^T x + \frac{1}{2} (m_1^T m_1 - m_2^T m_2)$$

posterior  $\leftarrow$   $P(C_1|x)$   $\leftarrow$  data likelihood  $\leftarrow$   $P(C_2|x)$   $\leftarrow$  prior

$$P(C_1|x) = \frac{P(x|C_1) P(C_1)}{P(x)} \quad \text{Bayes Rule}$$

$$P(C_2|x) = \frac{P(x|C_2) P(C_2)}{P(x)}$$

MAP rule

$\downarrow$   
Maximum a posteriori probability:  $\arg \max_i P(C_i|x)$

Gaussian Model:  $P(x|C_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} e^{-\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)}$

$$\arg \max_i \log p(C_i | x) = \arg \max_i \log(p(x | C_i) p(C_i))$$

$$\arg \max_i \boxed{\log p(x | C_i) + \log p(C_i)}$$

$$\boxed{\log p(x | C_i)} = -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)$$

$$\text{Decision Surface: } -\log p(C_1 | x) = -\log p(C_2 | x)$$

$$\boxed{\begin{aligned} \frac{d}{2} \log 2\pi + \frac{1}{2} \log |\Sigma_1| + \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) - \log p(C_1) = \\ \frac{d}{2} \log 2\pi + \frac{1}{2} \log |\Sigma_2| + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) - \log p(C_2) \end{aligned}}$$

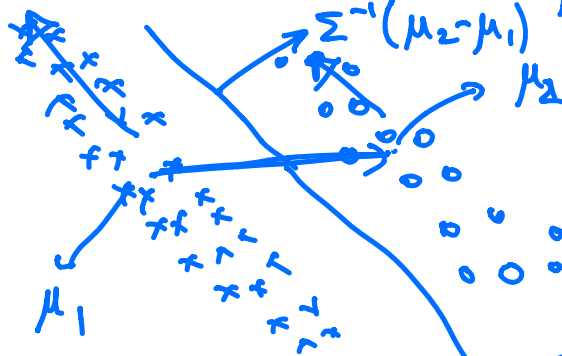
$$\text{Case I: } \Sigma_1 = \Sigma_2 = I$$

$$\frac{1}{2} (x - \mu_1)^T (x - \mu_1) - \log p(C_1) = \frac{1}{2} (x - \mu_2)^T (x - \mu_2) - \log p(C_2)$$

$$\frac{1}{2} (\cancel{x^T x} - 2\mu_1^T x + \|\mu_1\|^2) - \log p(C_1) =$$

$$\frac{1}{2} (\cancel{x^T x} - 2\mu_2^T x + \|\mu_2\|^2) - \log p(C_2)$$

$$(\mu_2 - \mu_1)^T x + \frac{1}{2} (\|\mu_1\|^2 - \|\mu_2\|^2) - \frac{\log p(C_1)}{p(C_2)} = 0$$



not necessarily  
parallel to  
line connecting

Case II:  $\Sigma_1 = \Sigma_2 = \Sigma$

(eigenvector  
of  $\Sigma$ )

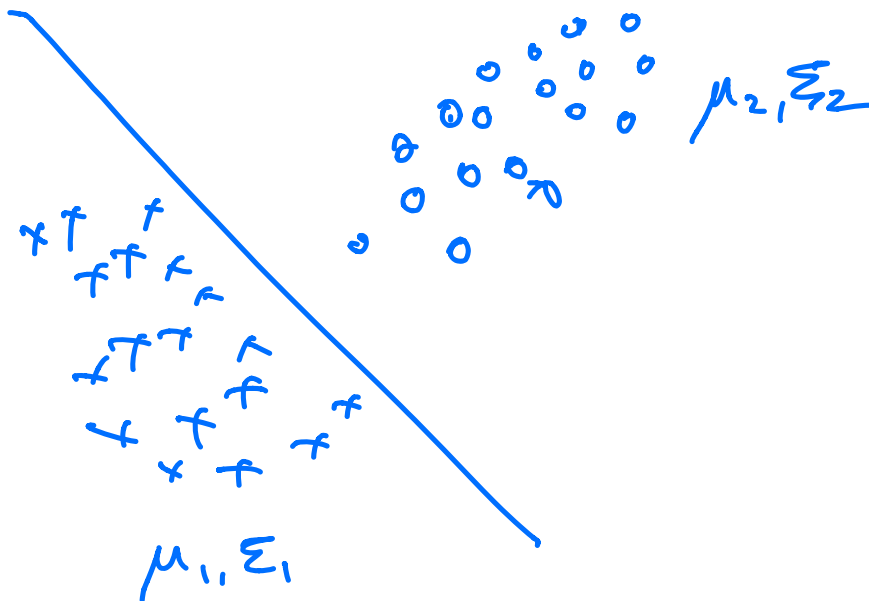
$$\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1) - \log p(c_1) = \frac{1}{2}(x-\mu_2)^T \Sigma^{-1}(x-\mu_2) - \log p(c_2)$$

Simplify:  $(\mu_2 - \mu_1)^T \Sigma^{-1} x + \frac{1}{2}(\mu_1 + \mu_2)^T \Sigma^{-1}(\mu_1 - \mu_2) - \log \frac{p(c_1)}{p(c_2)} = 0$

Decision Surface is Linear

$$w^T x + w_0 = 0$$

$$w = \Sigma^{-1}(\mu_2 - \mu_1)$$



# Classification: Regression Approaches

$$(x_i, y_i), \quad x_i \in \mathbb{R}^d, \quad y_i \in \mathbb{R}^k$$

$\nearrow$  k-th position  
 if  $x_i \in C_k$

$$y_i = [0, 0, \dots, 1, 0, \dots, 0]$$

$$w_0 + w_1 x(1) + w_2 x(2) + \dots + w_d x(d) \quad \text{Linear Fit}$$

Linear Discriminants:  $w_k^T x + w_{k_0}$  for k-th class

d+1 K (K-class problem)

$$\begin{array}{c}
 C_1 \\
 \vdots \\
 C_k
 \end{array}
 \left[ \begin{array}{c|cccc}
 1 & x_1(1) & x_1(2) & \dots & x_1(d) \\
 \hline
 \vdots & \vdots & \vdots & \ddots & \vdots \\
 \hline
 1 & x_N(1) & x_N(2) & \dots & x_N(d)
 \end{array} \right]
 \begin{array}{c}
 w_{10} \quad w_{20} \quad \dots \quad w_{k_0} \\
 \hline
 w_1 \quad w_2 \quad \dots \quad w_k
 \end{array}
 \begin{array}{c}
 y_1^T \\
 \vdots \\
 y_N^T
 \end{array}
 \left[ \begin{array}{cccc|cccc}
 1 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\
 \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\
 0 & 1 & 0 & \dots & 0 & 0 & \dots & 0 \\
 \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
 0 & 0 & 0 & \dots & 1 & 0 & \dots & 0 \\
 \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\
 0 & 0 & 0 & \dots & 0 & 0 & \dots & 1
 \end{array}
 \right]$$

$$\begin{aligned}
 XW &\approx Y \Rightarrow \min_W \frac{1}{2} \|XW - Y\|_F^2 \\
 &\equiv \min_W \frac{1}{2} \text{Tr}[(XW - Y)^T (XW - Y)]
 \end{aligned}$$

$$(X^T X) W^* = X^T Y \Rightarrow W^* = (X^T X)^{-1} X^T Y$$

$XW^*$  is prediction on training data

$x_i^T W$  could be  $\underline{[-5 \quad -3 \quad 0 \quad 3 \quad 2 \quad 1 \quad 0.5 \quad 0.1]}$

and this is not a good approximation to  
1-hot vectors

Least squares fit has obvious drawbacks

Special case of least squares fit:

2-class problem

$$y_i = \frac{N}{n_1} \text{ for each } x_i \in C_1, \quad n_1 = |C_1|$$

$$y_i = -\frac{N}{n_2} \text{ for each } x_i \in C_2$$

Turns out to be equivalent to:

Fisher's Linear Discriminant

### Logistic Regression

We had modeled each class as a Gaussian  
with covariance  $\Sigma$ :

$$\log \frac{p(C_i|x)}{p(C_j|x)} = \underbrace{\log \frac{p(C_i)}{p(C_j)}}_{w_0 + w^T x} - \frac{1}{2} (m_i + m_j)^T \Sigma^{-1} (m_i - m_j) + \underbrace{x^T \Sigma^{-1} (m_i - m_j)}_{\leftarrow}$$

k-class problem

$$\log \frac{p(c_1|x)}{p(c_k|x)} = w_{10} + w_1^T x \quad - \textcircled{1}$$

$$\log \frac{p(c_2|x)}{p(c_k|x)} = w_{20} + w_2^T x \quad - \textcircled{2}$$

⋮

$$\log \frac{p(c_{k-1}|x)}{p(c_k|x)} = w_{k-1,0} + w_{k-1}^T x \quad - \textcircled{k-1}$$

Write  $p_i = p(c_i|x)$

Add  $\textcircled{1}$ ,  $\textcircled{2}$ , ...,  $\textcircled{k-1}$

$$\textcircled{1} \Rightarrow \frac{p_1}{p_k} = e^{w_{10} + w_1^T x}$$

$$\frac{p_1}{p_k} + \frac{p_2}{p_k} + \dots + \frac{p_{k-1}}{p_k} = \sum_{i=1}^{k-1} e^{w_i^T x}$$

$$1 - p_k = p_k \sum_{i=1}^{k-1} e^{w_i^T x}$$

$$\Rightarrow p_k = \frac{1}{1 + \sum_{i=1}^{k-1} e^{w_i^T x}}$$

$$p_i = \frac{e^{w_i^T x}}{1 + \sum_{i=1}^{k-1} e^{w_i^T x}}, \quad i=1, 2, \dots, k-1$$



$\sigma(z) = \frac{1}{1+e^{-z}}$  is called the logistic

Sigmoid function

$$\sigma'(z) = \frac{1 \cdot e^{-z}}{(1+e^{-z})^2}$$

$$= \frac{e^{-z}}{1+e^{-z}} \cdot \frac{1}{1+e^{-z}}$$

$$= (1-\sigma(z)) \cdot \sigma(z)$$

$$\sigma' = \sigma(1-\sigma)$$

Parameters in logistic regression  $\{w_{i0}, w_i\}_{i=0}^{k-1}$  usually fit by maximum likelihood using the conditional likelihood given  $x$ .

Consider 2-class problem

$$\rightarrow p(C_1|x) = \boxed{p(w) = p}$$

$$\rightarrow p(C_2|x) = 1-p$$

$(x_i, y_i)$  is training data  
 $i=1, \dots, N$

Data Likelihood

$$\prod_{i=1}^N p^{y_i} (1-p)^{1-y_i}$$

$p(C_1|x)$  when  $x \in C_1$

let  $y_i = 1$  when  $x_i \in C_1$

$y_i = 0$  when  $x_i \in C_2$

Max log-likelihood

$$\max_w \ell(w) = \max_w \sum_{i=1}^N y_i \log p + (1-y_i) \log(1-p)$$

$$1-p = \frac{1}{1+e^{w^T x_i}}, p = \frac{e^{w^T x_i}}{1+e^{w^T x_i}} \Rightarrow \log p = w^T x_i - \log(1+e^{w^T x_i})$$

$\nabla_w \ell(w) = 0$  - Minimizer will satisfy

$$\nabla_w \log p = x_i - \frac{1}{1+e^{w^T x_i}} \cdot e^{w^T x_i} \cdot x_i$$

$$\nabla_w \log(1-p) = - \frac{e^{w^T x_i}}{1+e^{w^T x_i}} \cdot x_i$$

$$\nabla_w \ell(w) = \sum_{i=1}^N \left[ y_i x_i \left( 1 - \frac{e^{w^T x_i}}{1+e^{w^T x_i}} \right) \right. \\ \left. + (1-y_i) x_i \left( - \frac{e^{w^T x_i}}{1+e^{w^T x_i}} \right) \right]$$

$$= \sum x_i \left[ \cancel{y_i} - \frac{y_i e^{w^T x_i}}{1+e^{w^T x_i}} - \frac{e^{w^T x_i}}{1+e^{w^T x_i}} + \cancel{\frac{y_i e^{w^T x_i}}{1+e^{w^T x_i}}} \right]$$

$$= \sum_{i=1}^N x_i \left[ y_i - \frac{e^{w^T x_i}}{1+e^{w^T x_i}} \right] = 0$$

logistic regression parameters satisfy this

$$p(C_1 | x_i)$$

$$= \sum_{i=1}^N x_i \left[ y_i - (1 - \sigma(w^T x_i)) \right]$$

d+1 non-linear equations in  $w$

$\downarrow$   
(d+1) parameters

Re-weighted  
Iterated Least squares method (IRLS)

~~\*~~ Gradient Descent/Ascent

↙  
Newton's Method

↘  
Drawback is that each step  
of Gradient Descent is  $O(N)$

Stochastic Gradient Descent

Regularization

$$\lambda \|w\|_2^2$$

$$\lambda \|w\|_1$$